

Taller: **Data Quality Gate - Scripts**  
**Proyecto Big Data I**

*VS Code + conda + Python*

José Luis Gómez Ortega

2025

## 1 Intro (motivación)

Imagina que cargas encuestas de turismo (Alcoi) a tu base de datos para un dashboard. Algunos registros vienen con tildes en columnas, enteros como texto, combinaciones imposibles (*0 adultos y 0 menores*), valores ambiguos como “No entiendo la pregunta” y duplicados por marca temporal. Si eso llega a producción, pierdes tiempo y credibilidad.

Un **Data Quality Gate** es un guardia de frontera: **lee** tu dataset, **aplica reglas explícitas**, **genera un informe** y si detecta problemas críticos, **termina con código 1** para *bloquear la ingestión*. Hoy construimos ese gate como un **script** (no notebook): simple, reproducible y listo para enganchar en cualquier pipeline.

## 2 Objetivos

- Diferenciar **notebook** vs **script** y por qué el script es ideal para checks previos a ETL/ML.
- Ejecutar un script con **argumentos** y comprender el **exit code** (0/1).
- Dominar **rutas** con **pathlib** para evitar “file not found” según desde dónde ejecutes.
- Aplicar reglas de calidad: **esquema, tipos, dominios, rangos, unicidad y condicionales**.
- Generar un **informe CSV** de incidencias y usarlo para decidir *OK/Bloquear*.
- Entender cómo esto evita **DBs contaminadas, KPIs erróneos y modelos rotos**.

## 3 Conceptos clave (qué, dónde y por qué)

### Script vs Notebook

**Qué:** un script hace una tarea y termina (no interactivo). **Dónde:** pasos automatizados de pre-ingesta/pre-train. **Por qué:** reproducible, integrable (cron, Airflow, CI).

### CLI y argumentos

**Qué:** invocar `dq_check.py` con `-i`, `-r`, `-o`. **Dónde:** capa de orquestación del check. **Por qué:** configuración sin tocar código.

### Rutas con pathlib

**Qué:** construir rutas relativas al archivo (`__file__`) y no al *cwd*. **Dónde:** lectura de `data/`, `config/` y escritura en `out/`. **Por qué:** portabilidad real.

### Reglas de calidad

**Esquema, Tipos, Dominios, Rangos, Unicidad, Condicionales** (si origen=España entonces `ciudad_es` no nula; si Internacional entonces `pais_intl` no nulo; `n_adultos + n_menores > 1`).

### Informe & Exit code

**Qué:** CSV con incidencias + `exit(0/1)`. **Dónde:** artefacto en `out/`. **Por qué:** trazabilidad y automatización.

## 4 Flujo completo (todo resuelto y con archivos indicados)

### 4.1 Estructura mínima del proyecto

Listing 1: Árbol de directorios

```
1 project/  
2   data/  
3     raw/alcoy_sample.csv  
4   config/  
5     rules.yml  
6   out/  
7   scripts/  
8     dq_check.py  
9   requirements.txt  
10  README.md
```

### 4.2 Dependencias

**Archivo:** requirements.txt

Listing 2: requirements.txt

```
1 pandas  
2 pyyaml
```

### 4.3 Reglas preparadas (ajusta encabezados exactos si difieren)

**Archivo:** config/rules.yml

Listing 3: config/rules.yml

```
1 columns_map:  
2   timestamp: "Marca temporal"  
3   origen: "Cul es su procedencia?"  
4   ciudad_es: "En el caso de Espaa, qu ciudad?"  
5   pais_intl: "En el caso de Internacional, qu pas?"  
6   n_adultos: "# Con cuntas personas viajas? (adultos)"  
7   n_menores: "# Con cuntos menores viajas?"  
8   fin_semana: "Con el fin de semana"  
9  
10 schema:  
11   required: [timestamp, origen, n_adultos, n_menores, fin_semana]  
12  
13 types:  
14   timestamp: datetime  
15   origen: category  
16   ciudad_es: string  
17   pais_intl: string  
18   n_adultos: int  
19   n_menores: int  
20   fin_semana: category  
21  
22 domains:  
23   origen:  
24     allowed: ["Alcoi", "Espaa", "Internacional", "Comunidad Valenciana"]  
25   fin_semana:  
26     allowed: ["S", "No", "No entiendo la pregunta"]  
27  
28 ranges:
```

```

29     n_adultos: {min: 0, max: 10}
30     n_menores: {min: 0, max: 10}
31
32     uniques:
33         - [timestamp, origen]
34
35     conditionals:
36         - if: {origen: "Espaa"}
37           then: {require_not_null: ["ciudad_es"]}
38         - if: {origen: "Internacional"}
39           then: {require_not_null: ["pais_intl"]}
40         - rule: "n_adultos + n_menores >= 1"
41
42     severity:
43         missing_column: error
44         type_mismatch: error
45         domain_violation: warning
46         range_violation: warning
47         uniqueness_violation: error
48         conditional_violation: error

```

#### 4.4 Script principal (súper comentado)

Archivo: scripts/dq\_check.py

Listing 4: scripts/dq\_check.py

```

1  #!/usr/bin/env python3
2  """
3  Data Quality Gate simple para CSVs de encuestas (Alcoi).
4  - Lee CSV y reglas YAML.
5  - Mapea nombres reales -> cannicos.
6  - Fuerza tipos.
7  - Ejecuta checks y guarda informe CSV.
8  - Exit code: 0 (OK) / 1 (errores crticos).
9  """
10 import sys, argparse, logging
11 from pathlib import Path
12 import pandas as pd, yaml
13
14 # --- Rutas robustas relativas a este archivo ---
15 HERE = Path(__file__).resolve().parent      # ../project/scripts
16 PROJECT_ROOT = HERE.parent                  # ../project
17 DATA_DIR = PROJECT_ROOT / "data"
18 CONFIG_DIR = PROJECT_ROOT / "config"
19 OUT_DIR = PROJECT_ROOT / "out"
20
21 # --- Logging y argumentos CLI ---
22 def setup_logging(level: str = "INFO"):
23     logging.basicConfig(
24         level=getattr(logging, level.upper(), logging.INFO),
25         format="%(asctime)s | %(levelname)s | %(message)s",
26     )
27
28 def parse_args():
29     p = argparse.ArgumentParser(description="Data Quality Gate (simple)")
30     p.add_argument("-i", "--input", required=True, help="Ruta al CSV de entrada")
31     p.add_argument("-r", "--rules", default=str(CONFIG_DIR/"rules.yml"),
32                     help="Ruta al archivo de reglas YAML")
33     p.add_argument("-o", "--out", default=str(OUT_DIR),

```

```

34         help="Carpeta para el informe")
35     p.add_argument("--log-level", default="INFO")
36     return p.parse_args()
37
38 # --- Reglas (YAML) ---
39 def load_rules(path: Path) -> dict:
40     with open(path, "r", encoding="utf-8") as f:
41         data = yaml.safe_load(f)
42     return data or {}
43
44 # --- Renombrado y tipos ---
45 def rename_columns_to_canonical(df: pd.DataFrame, columns_map: dict) -> pd.DataFrame:
46     if not columns_map:
47         return df
48     real_to_canon = {v: k for k, v in columns_map.items()}
49     return df.rename(columns=lambda c: real_to_canon.get(c, c))
50
51 def coerce_types(df: pd.DataFrame, types_cfg: dict) -> pd.DataFrame:
52     if not types_cfg:
53         return df.copy()
54     out = df.copy()
55     for col, t in types_cfg.items():
56         if col not in out.columns: # si la columna no existe, seguimos
57             continue
58         try:
59             if t == "int":
60                 out[col] = pd.to_numeric(out[col], errors="coerce").astype("Int64")
61             elif t == "float":
62                 out[col] = pd.to_numeric(out[col], errors="coerce")
63             elif t == "datetime":
64                 out[col] = pd.to_datetime(out[col], errors="coerce", dayfirst=True)
65             elif t == "category":
66                 out[col] = out[col].astype("string").astype("category")
67             elif t == "string":
68                 out[col] = out[col].astype("string")
69         except Exception as e:
70             logging.warning("No se pudo convertir %s a %s: %s", col, t, e)
71     return out
72
73 # --- Checks: devuelven lista de 'findings' dicts ---
74 def check_schema(df, rules):
75     required = set((rules.get("schema") or {}).get("required", []))
76     missing = sorted(list(required - set(df.columns)))
77     findings = []
78     for col in missing:
79         findings.append({"rule": "missing_column", "column": col, "row": None,
80                         "detail": "Columna obligatoria ausente"})
81     return findings
82
83 def check_domains(df, rules):
84     findings = []
85     for col, cfg in (rules.get("domains") or {}).items():
86         if col not in df.columns:
87             continue
88         allowed = set(cfg.get("allowed", []))
89         mask = (~df[col].isna()) & (~df[col].isin(allowed))
90         for idx in df[mask].index:
91             findings.append({"rule": "domain_violation", "column": col, "row": int(idx),
92                             "detail": f"Valor={df.at[idx,col]!r} no permitido"})
93     return findings

```

```

94
95 def check_ranges(df, rules):
96     findings = []
97     for col, cfg in (rules.get("ranges") or {}).items():
98         if col not in df.columns:
99             continue
100         mn, mx = cfg.get("min"), cfg.get("max")
101         if mn is not None:
102             bad = (df[col] < mn)
103             for idx in df[bad.fillna(False)].index:
104                 findings.append({"rule": "range_violation", "column": col, "row": int(idx),
105                                "detail": f"< min({mn})"})
106         if mx is not None:
107             bad = (df[col] > mx)
108             for idx in df[bad.fillna(False)].index:
109                 findings.append({"rule": "range_violation", "column": col, "row": int(idx),
110                                "detail": f"> max({mx})"})
111     return findings
112
113 def check_uniques(df, rules):
114     findings = []
115     for cols in (rules.get("uniques") or []):
116         dups = df.duplicated(subset=cols, keep=False)
117         for idx in df[dups].index:
118             findings.append({"rule": "uniqueness_violation",
119                             "column": ", ".join(cols), "row": int(idx),
120                             "detail": "Fila duplicada para la clave dada"})
121     return findings
122
123 def check_conditionals(df, rules):
124     findings = []
125     for cond in (rules.get("conditionals") or []):
126         if "rule" in cond:
127             expr = cond["rule"] # p.ej. "n_adultos + n_menores >= 1"
128             bad = ~df.eval(expr) # incumplimientos
129             for idx in df[bad.fillna(True)].index:
130                 findings.append({"rule": "conditional_violation", "column": None,
131                                "row": int(idx), "detail": f"No cumple: {expr}"})
132         else:
133             if_ = cond.get("if", {})
134             then = cond.get("then", {})
135             mask = pd.Series(True, index=df.index)
136             for c, v in if_.items():
137                 mask = mask & (df[c] == v)
138             req = then.get("require_not_null", []) or []
139             for col in req:
140                 bad = mask & df[col].isna()
141                 for idx in df[bad].index:
142                     findings.append({"rule": "conditional_violation", "column": col,
143                                    "row": int(idx), "detail": f"{col} requerido cuando {if_}"})
144     return findings
145
146 def apply_severity(findings, rules):
147     sev_map = (rules.get("severity") or {})
148     for f in findings:
149         f["severity"] = sev_map.get(f["rule"], "warning")
150
151 # --- Programa principal ---
152 def main():
153     args = parse_args()

```

```

154     setup_logging(args.log_level)
155
156     out_dir = Path(args.out)
157     out_dir.mkdir(parents=True, exist_ok=True)
158
159     logging.info("Leyendo CSV: %s", args.input)
160     df_raw = pd.read_csv(args.input) # ajusta sep/encoding si tu CSV lo requiere
161     rules = load_rules(Path(args.rules))
162
163     df = rename_columns_to_canonical(df_raw, rules.get("columns_map"))
164     df = coerce_types(df, rules.get("types"))
165
166     findings = []
167     findings += check_schema(df, rules)
168     findings += check_domains(df, rules)
169     findings += check_ranges(df, rules)
170     findings += check_uniques(df, rules)
171     findings += check_conditionals(df, rules)
172
173     apply_severity(findings, rules)
174     summary = {}
175     for f in findings:
176         summary[f["severity"]] = summary.get(f["severity"], 0) + 1
177     logging.info("Incidencias por severidad: %s", summary or "sin incidencias")
178
179     report_path = out_dir / "dq_report.csv"
180     pd.DataFrame(findings).to_csv(report_path, index=False)
181     logging.info("Informe: %s", report_path)
182
183     has_errors = any(f["severity"] == "error" for f in findings)
184     if has_errors:
185         logging.error("Errores criticos detectados. Abortando (exit=1).")
186         sys.exit(1)
187     else:
188         logging.info("Calidad OK. (exit=0)")
189         sys.exit(0)
190
191 if __name__ == "__main__":
192     main()

```

## 4.5 Ejecución base (con salida y exit code)

### Terminal (Linux/Mac):

Listing 5: Comando base

```

1 python scripts/dq_check.py -i data/raw/alcoy_sample.csv -r config/rules.yml -o out
2 echo $?      # muestra 0 si OK, 1 si errores criticos

```

### PowerShell (Windows):

```

1 python scripts/dq_check.py -i data\raw\alcoy_sample.csv -r config\rules.yml -o out
2 $LASTEXITCODE

```

**Resultado esperado:** se genera out/dq\_report.csv. En consola se imprime un resumen por severidad. Si hay severidad error, el exit code es 1 (bloquearías la carga a BBDD).

## 4.6 Forzar un fallo y entender la política de severidad

**Archivo:** config/rules.yml

Cambia temporalmente el rango para provocar incidencias y observar el *exit code*:

```
1 ranges:
2   n_adultos: {min: 5, max: 5}  # absurdo a proposito (slo 5)
```

Ejecuta de nuevo y mira out/dq\_report.csv. Si sólo aparecen warning, eleva la severidad:

```
1 severity:
2   range_violation: error
```

Vuelve a ejecutar: ahora el *exit code* debe ser 1.

## 4.7 Condicionales típicas del dataset de Alcoi

Ya incluidas en config/rules.yml:

- Si origen=España  $\Rightarrow$  ciudad\_es no nula.
- Si origen=Internacional  $\Rightarrow$  pais\_intl no nulo.
- `n_adultos + n_menores > 1`.

Edita 1–2 filas del CSV para simular incumplimientos y comprueba que aparecen como `conditional_violation` (por defecto, `error`).

## 4.8 Duplicados (unicidad)

En config/rules.yml:

```
1 uniques:
2   - [timestamp, origen]
```

Duplica una fila con la misma pareja y ejecuta: verás `uniqueness_violation` y *exit* = 1.

## 4.9 Gotchas típicos & soluciones

- **Rutas:** construimos con `Path(__file__)`; no depende del *cwd*.
- **Encoding/Separador:** si tu CSV es *latin-1* o usa `'`;', edita en `scripts/dq_check.py`, función `main()`, línea de lectura:

```
1 df_raw = pd.read_csv(args.input, sep=";", encoding="latin-1")
```
- **Tildes y espacios en encabezados:** por eso mapeamos a nombres canónicos (`config/rules.yml`  $\rightarrow$  `columns_map`).

## 5 Guía exprés para el/la docente

1. Prepara el árbol de proyecto y coloca tu CSV real en `data/raw/alcoy_sample.csv`.
2. Ajusta `config/rules.yml`  $\rightarrow$  `columns_map` a los encabezados *exactos*.
3. Instala dependencias: `pip install -r requirements.txt`.
4. Ejecuta el script y abre `out/dq_report.csv` en clase.



5. Muestra cómo cambiar **severidades** para forzar *exit* = 1.
6. Conecta con casos reales: “*esto corre antes de to\_sql()*”.

## 6 FAQ

**No encuentra rules.yml:** ejecuta desde la raíz del proyecto.

**UnicodeDecodeError:** añade `encoding="latin-1"` en `pd.read_csv`.

**ParserError:** añade `sep=";"` en `pd.read_csv`.

**No salen incidencias:** revisa `columns_map` (debe coincidir *exacto* con los encabezados del CSV).

*Un script sencillo hoy = horas de limpieza ahorradas mañana.*

**Data Quality Gate** antes de cualquier ETL/ML/BI. #shipit